

Invenio@HGF – Collaborative repository infrastructure

Open Repositories 2014 - Helsinki

Alexander Wagner¹, Robert Thiele²
for the Collaboration

¹Forschungszentrum Jülich, Zentralbibliothek

²DESY Hamburg, Bibliothek

13. June 2014





- Partner
- Initial TODO
- Accomplishments
- Lessons learned
- Project group

Project Partners



Deutsches Elektronensynchrotron, Zentralbibliothek

≈ 2000 + 3000



Forschungszentrum Jülich, Zentralbibliothek

≈ 5000 + 1000



GSI Helmholtzzentrum für Schwerionenforschung, Bibliothek + Base-IT

≈ 1050



Deutsches Krebsforschungszentrum, Bibliothek

≈ 3000



Maier-Leibniz-Zentrum, Garching

≈ 300



RWTH Aachen, Hochschulbibliothek

≈ 9000

Museum Zitadelle Jülich



Institut für Experimentelle Kernphysik, Karlsruhe



Project Partners



Deutsches Elektronensynchrotron, Zentralbibliothek

≈ 2000 + 3000



Forschungszentrum Jülich, Zentralbibliothek

≈ 5000 + 1000



GSI Helmholtzzentrum für Schwerionenforschung, Bibliothek + Base-IT

≈ 1050



Deutsches Krebsforschungszentrum, Bibliothek

≈ 3000



Maier-Leibniz-Zentrum, Garching

≈ 300



RWTH Aachen, Hochschulbibliothek

≈ 9000

Museum Zitadelle Jülich



Institut für Experimentelle Kernphysik, Karlsruhe

Open for new Partners!



Project Partners



Deutsches Elektronensynchrotron, Zentralbibliothek

≈ 2000 + 3000



Forschungszentrum Jülich, Zentralbibliothek

≈ 5000 + 1000



GSI Helmholtzzentrum für Schwerionenforschung, Bibliothek + Base-IT

≈ 1050



Deutsches Krebsforschungszentrum, Bibliothek

≈ 3000



Maier-Leibniz-Zentrum, Garching

≈ 300



RWTH Aachen, Hochschulbibliothek

≈ 9000

Museum Zitadelle Jülich



Institut für Experimentelle Kernphysik, Karlsruhe

Open for new Partners!

Serving now ≈ **20.000** people (+ visitors)

(≈ 260.000 documents + 80.000 Authorities)





GERMAN
CANCER RESEARCH CENTER
IN THE HELMHOLTZ ASSOCIATION

50 Years – Research for
A Life Without Cancer

- Largest German bio-medicine research center
- 3000 people:
1000 researchers inside
- over 90 divisions and groups: research of cancer-producing, risk factors and strategies against cancer
- Replacing of own existing repository system with Invenio@HGF
- Roll out planned for mid of 2015



Goal

Replace existing systems, at GSI build up from scratch.
User-centric design (users aka scientists)



Goal

Replace existing systems, at GSI build up from scratch.
User-centric design (users aka scientists)

- 1 “Learn Invenio” (thanks to CERN ☺)



Goal

Replace existing systems, at GSI build up from scratch.
User-centric design (users aka scientists)

- 1 “Learn Invenio” (thanks to CERN ☺)
- 2 Define wording. . . (different institutions!)



Goal

Replace existing systems, at GSI build up from scratch.
User-centric design (users aka scientists)

- 1 “Learn Invenio” (thanks to CERN ☺)
- 2 Define wording. . . (different institutions!)
- 3 Build infrastructure: git and friends



Goal

Replace existing systems, at GSI build up from scratch.
User-centric design (users aka scientists)

- 1 “Learn Invenio” (thanks to CERN ☺)
- 2 Define wording. . . (different institutions!)
- 3 Build infrastructure: git and friends
- 4 Build more infrastructure: authorities and friends



Goal

Replace existing systems, at GSI build up from scratch.
User-centric design (users aka scientists)

- 1 “Learn Invenio” (thanks to CERN ☺)
- 2 Define wording. . . (different institutions!)
- 3 Build infrastructure: git and friends
- 4 Build more infrastructure: authorities and friends
- 5 Build a **deployment scheme**: InstallInvenio and friends



Goal

Replace existing systems, at GSI build up from scratch.
User-centric design (users aka scientists)

- 1 “Learn Invenio” (thanks to CERN ☺)
- 2 Define wording. . . (different institutions!)
- 3 Build infrastructure: git and friends
- 4 Build more infrastructure: authorities and friends
- 5 Build a **deployment scheme**: InstallInvenio and friends

We need to roll out 10+ instances



Goal

Replace existing systems, at GSI build up from scratch.
User-centric design (users aka scientists)

- 1 “Learn Invenio” (thanks to CERN ☺)
- 2 Define wording. . . (different institutions!)
- 3 Build infrastructure: git and friends
- 4 Build more infrastructure: authorities and friends
- 5 Build a **deployment scheme**: InstallInvenio and friends

We need to roll out 10+ instances
with different data sets



Goal

Replace existing systems, at GSI build up from scratch.
User-centric design (users aka scientists)

- 1 “Learn Invenio” (thanks to CERN ☺)
- 2 Define wording. . . (different institutions!)
- 3 Build infrastructure: git and friends
- 4 Build more infrastructure: authorities and friends
- 5 Build a **deployment scheme**: InstallInvenio and friends

We need to roll out 10+ instances
with different data sets and keep them consistent on code level

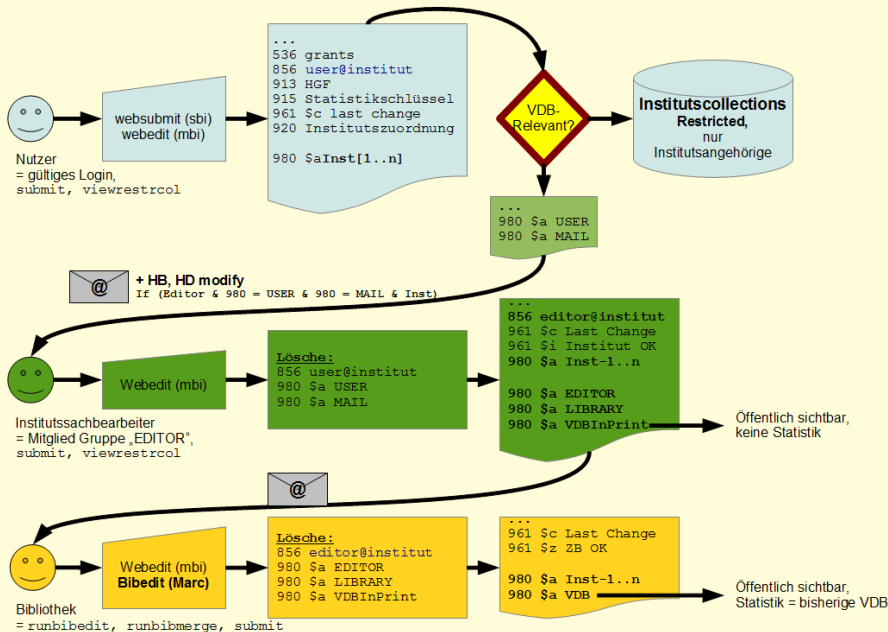


- Design the system around [web based literature management](#)



- Design the system around [web based literature management](#)
- Design a document workflow





- Design the system around [web based literature management](#)
- Design document workflow ([3 steps with privilege escalation](#))
- Design easy ingestion workflow ([websubmit, imports, author disambiguation](#))



Import data ⓘ DOI, arXiv, PUBMED...

ISBN ⓘ 978-3-642-12893-6

Group(s) involved * ⓘ

T: Theorie-Gruppe x

Select or type in name, shortcut (e.g. ATLAS, FS-PE, MKK)

Beamline/Experiment/Facility (Automatically assigns Grants) * ⓘ

514 - Theoretical Particle Physics (POF II: 2010 - 2014) x

Select PETRA beamline, HERA, facility machine,...

Grant name / Proposal No. ⓘ

514 - Theoretical Particle Physics (POF II: 2010 - 2014) x

EU project, FS proposal number (e.g. I-20120768)

Report Number ⓘ

DESY 14-075

Author(s) / Contributor(s) * ⓘ

Bonvini, Marco -> Bonvini, Marco (>DESY / T) Corresponding Author x

Marzani, Simone [Extern] Author ✓ x

Start typing lastname and select...

Title * ⓘ

Resummed Higgs cross section at $N^3\text{LL}$

Title preview:

Resummed Higgs cross section at $N^3\text{LL}$

Journal ⓘ [Type name, issn...; use "" for exact match or fields e.g. title:"nature"]

DOI ⓘ Use Import data for automatic prefill

Volume ⓘ **Issue** ⓘ **Pages** ⓘ e.g. 47-103

Publication Year * ⓘ 2014 **Language** ⓘ Click to select...

Edition ⓘ 5th ed. **Publisher** ⓘ Forschungszentrum Jülich, Verlag

Place of publication ⓘ Jülich

Title Series ⓘ

Abstract ⓘ

We present accurate predictions for the inclusive production of a Higgs boson in proton-proton collisions, via gluon-gluon fusion. Our calculation includes next-to-next-to-leading order (NNLO) corrections in perturbative QCD, as well as the resummation of threshold-enhanced contributions to next-to-next-to-leading logarithmic (NNLL) accuracy, with

- Design the system around **web based literature management**
- Design document workflow (3 steps with privilege escalation)
- Design easy ingestion workflow (websubmit, imports, author disambiguation)



- Design the system around **web based literature management**
- Design document workflow (3 steps with privilege escalation)
- Design easy ingestion workflow (websubmit, imports, author disambiguation)

Finally we wrote some code...

- Design the system around **web based literature management**
- Design document workflow (3 steps with privilege escalation)
- Design easy ingestion workflow (websubmit, imports, author disambiguation)

Finally we wrote some code...

Every unwritten line is a good line



Initial ToDo

- Design the system around **web based literature management**
- Design document workflow (3 steps with privilege escalation)
- Design easy ingestion workflow (websubmit, imports, author disambiguation)

Finally we wrote some code...

Every unwritten line is a good line, still: \approx 55.000 lines



- Design the system around **web based literature management**
- Design document workflow (3 steps with privilege escalation)
- Design easy ingestion workflow (websubmit, imports, author disambiguation)

Finally we wrote some code...

Every unwritten line is a good line, still: \approx 55.000 lines

- Migrate old data (various, proprietary sources)
- Train the inputters and users (secretaries, scientists, librarians)
- Hook up with content management system(s) (visibility!)



Content management system(s)

The screenshot shows the website of the Peter Grünberg Institut (PGI) at Jülich. The header includes the Jülich Forschungszentrum logo and navigation links. The main content area displays a list of publications under the heading "Referierte Zeitschriftenbeiträge 2013". The list includes entries such as "Stoner-Pauling behavior in LiMgPdsn-type multifunctional quarter Heusler materials" and "GW study of topological insulators".

The screenshot shows the Photon Science Publications list on the DESY website. The header includes the DESY logo and navigation links. The main content area displays a list of publications under the heading "Photon Science Publications in the DESY Publication Database (PUBDOB)". The list includes entries such as "Structural Differences Explain Diverse Functions of Photosystem Actin" and "Frontend preserving channel-cut optics for coherent x-ray scattering experiments".



- Design a document workflow (3 steps with privilege escalation)
- Establish easy ingestion workflow (websubmit, imports, author disambiguation)

Finally we wrote some code...

Every unwritten line is a good line, still: ≈ 55.000 lines

- Migrate old data (various, proprietary sources)
- Train the inputters and users (secretaries, scientists, librarians)
- Hook up with content management system(s) (visibility!)
- Derive necessary reporting (statistics for the Helmholtz Foundation etc.)

- Design a document workflow (3 steps with privilege escalation)
- Establish easy ingestion workflow (websubmit, imports, author disambiguation)

Finally we wrote some code. . .

Every unwritten line is a good line, still: ≈ 55.000 lines

- Migrate old data (various, proprietary sources)
- Train the inputters and users (secretaries, scientists, librarians)
- Hook up with content management system(s) (visibility!)
- Derive necessary reporting (statistics for the Helmholtz Foundation etc.)
- **Get it up and running** (First Light: 11/19/2012)



Accomplishments and status

The collage displays four distinct scientific databases and their search capabilities:

- RWTH Aachen University:** A screenshot of the university's official website, showing navigation links and institutional information.
- PUBDB (Particle and Astrophysics Database):** A search interface for particle and astrophysics data, featuring a search bar, filters, and a list of results.
- IMPULSE (Institute for Materials Physics):** A search interface for materials science data, including a search bar, filters, and a list of results.
- GSI Repository:** A search interface for GSI (Gesellschaft für Schwerionenforschung) data, featuring a search bar, filters, and a list of results.



- All partners have running systems (roll out works)



Accomplishments and status

- All partners have running systems (roll out works)
- Almost all partners are online



Accomplishments and status

- All partners have running systems (roll out works)
- Almost all partners are online
- Rich websubmit (including **repeatable field** handling)
- Importer routines (doi, pmid, arXiv, inspire, ISBN, own recs, . . . in **websubmit**)



Accomplishments and status

- All partners have running systems (roll out works)
- Almost all partners are online
- Rich websubmit (including **repeatable field** handling)
- Importer routines (doi, pmid, arXiv, inspire, ISBN, own recs, . . . in **websubmit**)
- **Authorities**



- All partners have running systems (roll out works)
- Almost all partners are online
- Rich websubmit (including **repeatable field** handling)
- Importer routines (doi, pmid, arXiv, inspire, ISBN, own recs,... in **websubmit**)
- **Authorities**
 - **Generate** (\approx 80.000 recs)
 - **Use** (e. g. JSON returns, statistics...)
 - **Share** (MarcXML OAI-PMH)



Accomplishments and status

- All partners have running systems (roll out works)
- Almost all partners are online
- Rich websubmit (including **repeatable field** handling)
- Importer routines (doi, pmid, arXiv, inspire, ISBN, own recs,... in **websubmit**)
- **Authorities**
 - Generate (\approx 80.000 recs)
 - Use (e. g. JSON returns, statistics...)
 - Share (MarcXML OAI-PMH)
- Implement



Accomplishments and status

- All partners have running systems (roll out works)
- Almost all partners are online
- Rich websubmit (including **repeatable field** handling)
- Importer routines (doi, pmid, arXiv, inspire, ISBN, own recs, ... in **websubmit**)
- **Authorities**
 - **Generate** (≈ 80.000 recs)
 - **Use** (e. g. JSON returns, statistics...)
 - **Share** (MarcXML OAI-PMH)
- **Implement**
 - **Author identification** (ORCID ready!)
 - **Output formats** (JSON, BibTeX, EndNote... or special formats for our partners)
 - **Reporting** (publication statistics)
 - **Delivery to content management systems**



Analyzing collection "VDB" for WEB year 2013

27216 records in collection "VDB"

2527 records relevant for "WEB 2013"

	records	JCR listed	(JCR)	Refused 2013 Ref
6	records	WOS and not JCR	(WOSnotJCR)	Refused 2013 Ref
777	records	Web of Science	(WOS)	Refused 2013 Ref
777	records	WOS listed journal OR entry	(WOS)	Refused 2013 Ref
Per program and statistics key				
	records	Other refereed	(Other)	Refused 2013 Ref
48	records	Refereed NOT in WoS	(Other)	Refused 2013 Ref
40	records	All refereed	(Ref.)	Refused 2013 Ref
1017	records	Pubmed listed	(Med)	Refused 2013 Ref
800	records	Scopus listed	(Scop)	Refused 2013 Ref
112	records	DOAJ listed	(DOAJ)	Refused 2013 Ref
2	of type	Diploma Thesis		Refused 2013 Ref
56	of type	Dissertation / PhD Thesis		Refused 2013 Ref
1	of type	Habil / Postdoctoral Thesis (Non-per)		Refused 2013 Ref
711	of type	Internal Report		Refused 2013 Ref
1102	of type	Journal Article		Refused 2013 Ref
2	of type	Lecture		Refused 2013 Ref
14	of type	Master Thesis		Refused 2013 Ref
8	of type	Bachelor Thesis		Refused 2013 Ref
213	of type	Poster		Refused 2013 Ref
...

- Workflow



■ Workflow

- **Webbaskets** (e. g. revision lists)
- **Alerts** (e. g. revision lists)
- **Collections** (e. g. private for institutes)
- **Webmessage** (e. g. correction requests)



- Workflow
 - Webbaskets (e. g. revision lists)
 - Alerts (e. g. revision lists)
 - Collections (e. g. private for institutes)
 - Webmessage (e. g. correction requests)
- Authority records (almost everywhere)



- **Workflow**
 - **Webbaskets** (e. g. revision lists)
 - **Alerts** (e. g. revision lists)
 - **Collections** (e. g. private for institutes)
 - **Webmessage** (e. g. correction requests)
- **Authority records** (almost everywhere)
- **OAI-PMH** (authority exchange)

- Workflow
 - Webbaskets (e. g. revision lists)
 - Alerts (e. g. revision lists)
 - Collections (e. g. private for institutes)
 - Webmessage (e. g. correction requests)
- Authority records (almost everywhere)
- OAI-PMH (authority exchange)
- High-level API (setup: e. g. collections, roles, groups, baskets... ; **no db-dump** sharing)



- Workflow
 - Webbaskets (e. g. revision lists)
 - Alerts (e. g. revision lists)
 - Collections (e. g. private for institutes)
 - Webmessage (e. g. correction requests)
- Authority records (almost everywhere)
- OAI-PMH (authority exchange)
- High-level API (setup: e. g. collections, roles, groups, baskets... ; **no db-dump** sharing)
- jQuery/jQueryUI (websubmit)



- Workflow
 - Webbaskets (e. g. revision lists)
 - Alerts (e. g. revision lists)
 - Collections (e. g. private for institutes)
 - Webmessage (e. g. correction requests)
- Authority records (almost everywhere)
- OAI-PMH (authority exchange)
- High-level API (setup: e. g. collections, roles, groups, baskets... ; **no db-dump** sharing)
- jQuery/jQueryUI (websubmit)
- intbitset (e. g. statistics)



- CERN is way to fast to keep up with



Lessons learned / Next steps

- CERN is **way to fast** to keep up with
- Never use Dublin Core again (complex migration, to few data fields...)



Lessons learned / Next steps

- CERN is **way to fast** to keep up with
- Never use Dublin Core again (complex migration, to few data fields...)
- All libraries are the same ☺



- CERN is **way to fast** to keep up with
- Never use Dublin Core again (complex migration, to few data fields...)
- All libraries are the same ☺
- **Upgrade to 1.1.x:**
 - get OAI-Server fixed, no hanging bibsched, etc.
 - testing and bugfixing on our test systems
 - roll out update in July by our partners...



- CERN is **way to fast** to keep up with
- Never use Dublin Core again (complex migration, to few data fields...)
- All libraries are the same ☺
- **Upgrade to 1.1.x:**
 - get OAI-Server fixed, no hanging bibsched, etc.
 - testing and bugfixing on our test systems
 - roll out update in July by our partners...

However...

In our use case switching of the base system is non-trivial

(Remember: 10+ instances...)



- CERN is **way to fast** to keep up with
- Never use Dublin Core again (complex migration, to few data fields...)
- All libraries are the same ☺
- **Upgrade to 1.1.x:**
 - get OAI-Server fixed, no hanging bibsched, etc.
 - testing and bugfixing on our test systems
 - roll out update in July by our partners...

However...

In our use case switching of the base system is non-trivial

(Remember: 10+ instances...)

- Open up for new partners
- Clean up our code and give it back → moving to github



- *Martin Köhler*^a
- *Robert Thiele*^a
- *Katrin Große*^b
- *Stefan Hesselbach*^c
- *Bernhard Mittermaier*^d
- *Anna Fründ*^d
- *Heike Lexis*^d
- *Cornelia Plott*^d
- *Christoph Holzke*^d
- *Alexander Wagner*^d

- *Dagmar Sitek*^e
- *Gudrun Friedburg*^e
- *Jürgen Neuhaus*^f
- *Connie Hesse*^f
- *Björn Pedersen*^f
- *Ulrike Eich*^g
- *Louai Barake*^g
- *Abdoulaye Diallo*^g
- *Roland Rappmann*^g
- *Dominik Schmitz*^g
- *Edmund Wollgarten*^g

^a DESY Library and Documentation; ^b GSI Library and Documentation; ^c GSI Base-IT;

^d Forschungszentrum Jülich, Zentralbibliothek; ^e DKFZ Heidelberg; ^f MLZ, Garching; ^g RWTH Aachen, Hochschulbibliothek

- **Invenio @ HGF - Technical background**

[Talk at Invenio Developer Forum](#)

- **Collaborative tools for an institutional repository**

[Talk at Helmholtz OA Workshop](#)

- **JuSER – Publications Database**

[Introductory course at Jülich](#)

- **JuSER - Autorenhandling**

[Talk at HGF-ORCID Meeting, Berlin \(in german\)](#)

- **Invenio @HGF – status and perspectives**

[Talk at 2nd Invenio User Group Workshop, \[sic!\], Jülich, Germany](#)

- **The Helmholtz INVENIO Repository Project**

[Talk at SACITIL-2014, Kolkata, India](#)



Thanks!



Robert Thiele
DESY-Bibliothek

Subject Specialist for Photon
Science

Tel.: +49-40-8998-1927
robert.thiele@desy.de

This document is available as

DESY-2014-02793 or FZJ-2014-02848



Typeset by pdfL^AT_EX

